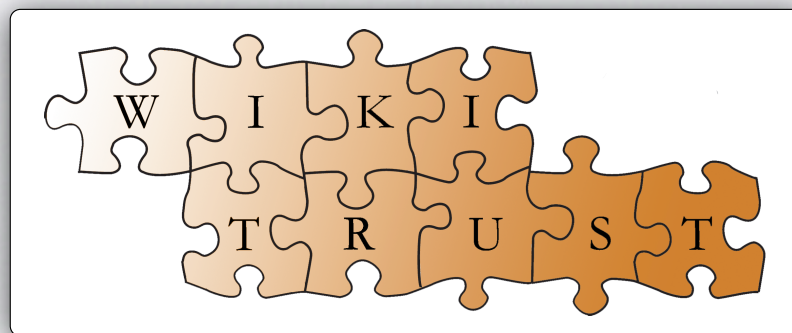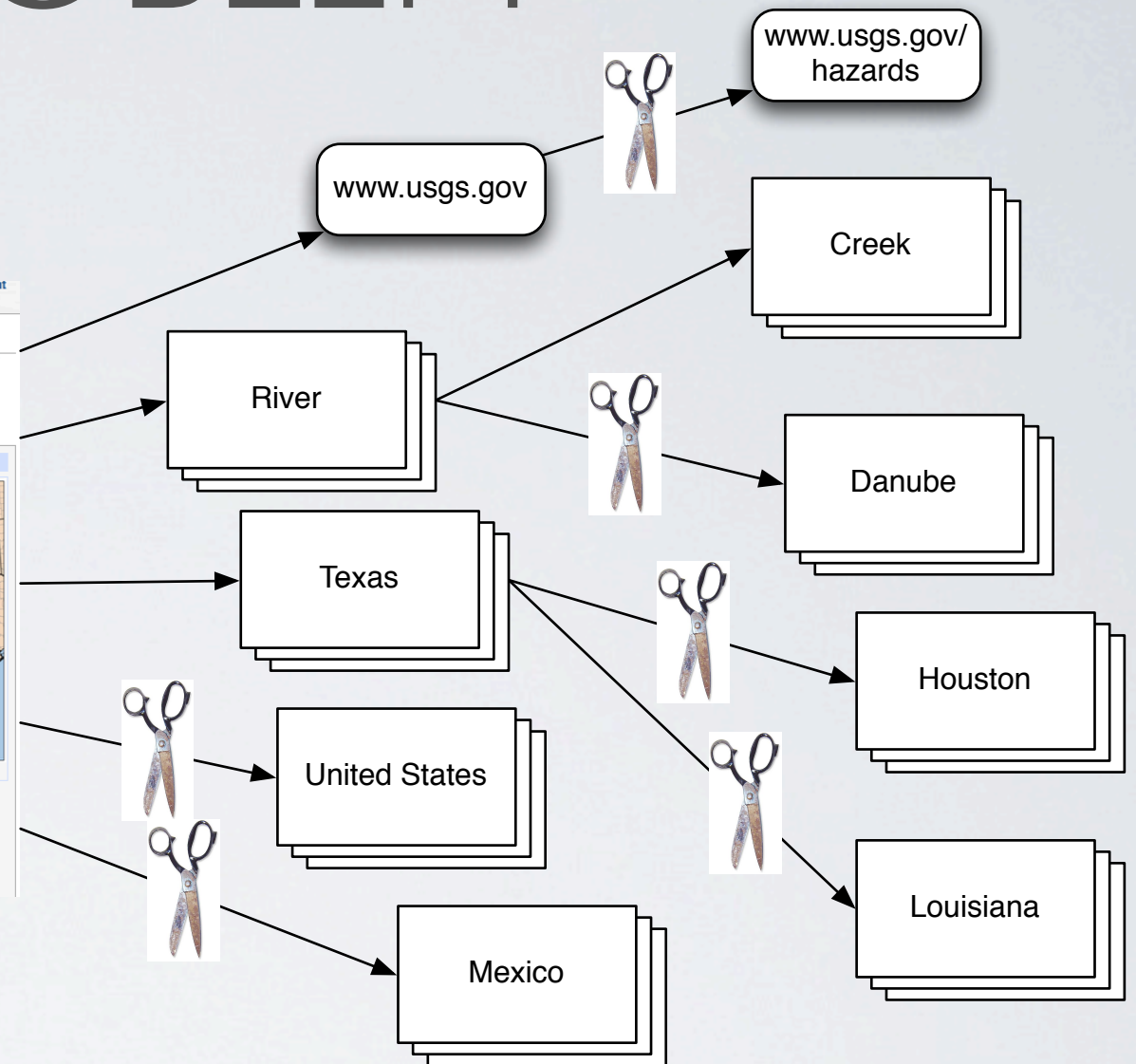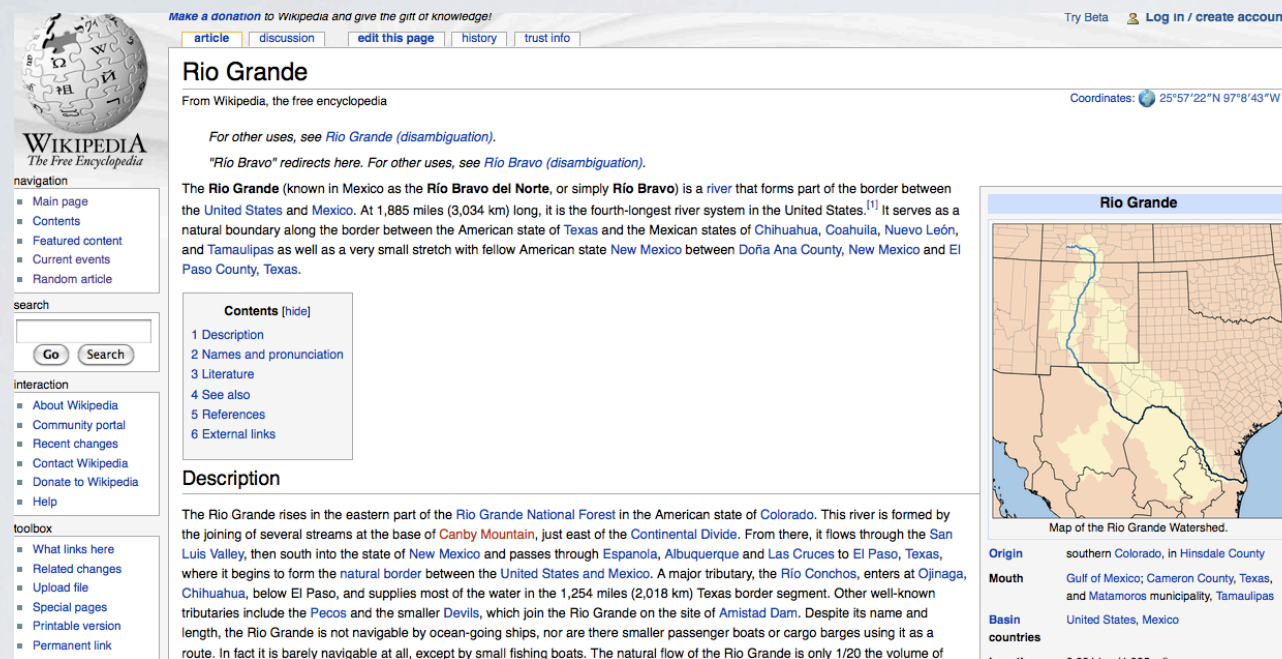# SCRAWL: A SEMANTIC CRAWLER FOR THE WIKIPEDIA AND BEYOND

Ian Pye, Luca de Alfaro
Shelly Spearing, Jorge Roman

# THE PROBLEM



The Wikipedia is too big! We just want to look at the (potentially) interesting parts.

# OUR SOLUTION

Target Page

Page to Evaluate

HTML

orch8.net, extract tri-grams

Target Histogram

Close Enough?

No

Yes

Disk

Revision Histogram

# OUR SOLUTION

## 1) Download and Render to HTML

Target Page

Page to Evaluate

HTML

# OUR SOLUTION
## 1) Download and Render to HTML
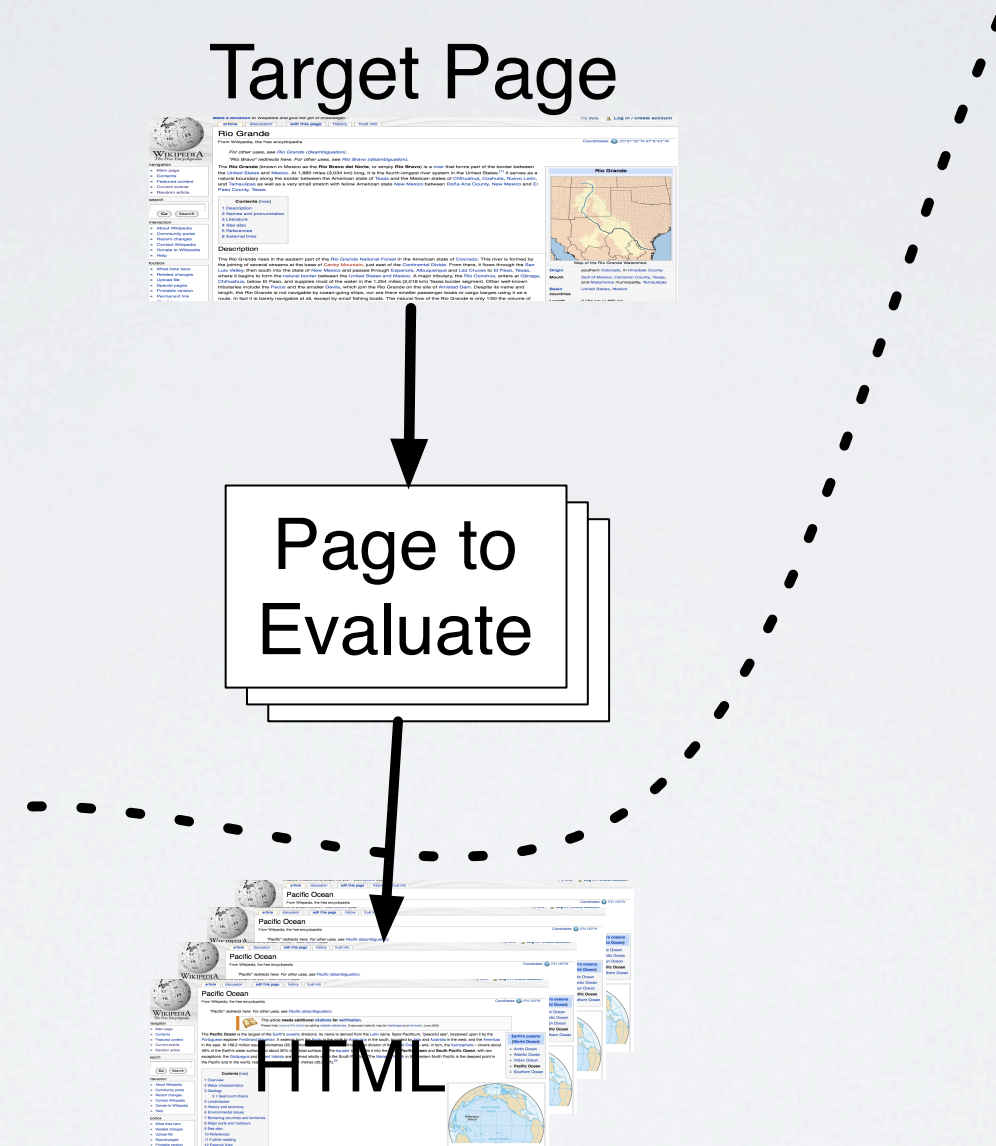## 2) From HTML to Semantic Digest

Page to
Evaluate

HTML

orch8.net,
extract
tri-grams

Revision
Histogram

# OUR SOLUTION
## 1) Download and Render to HTML
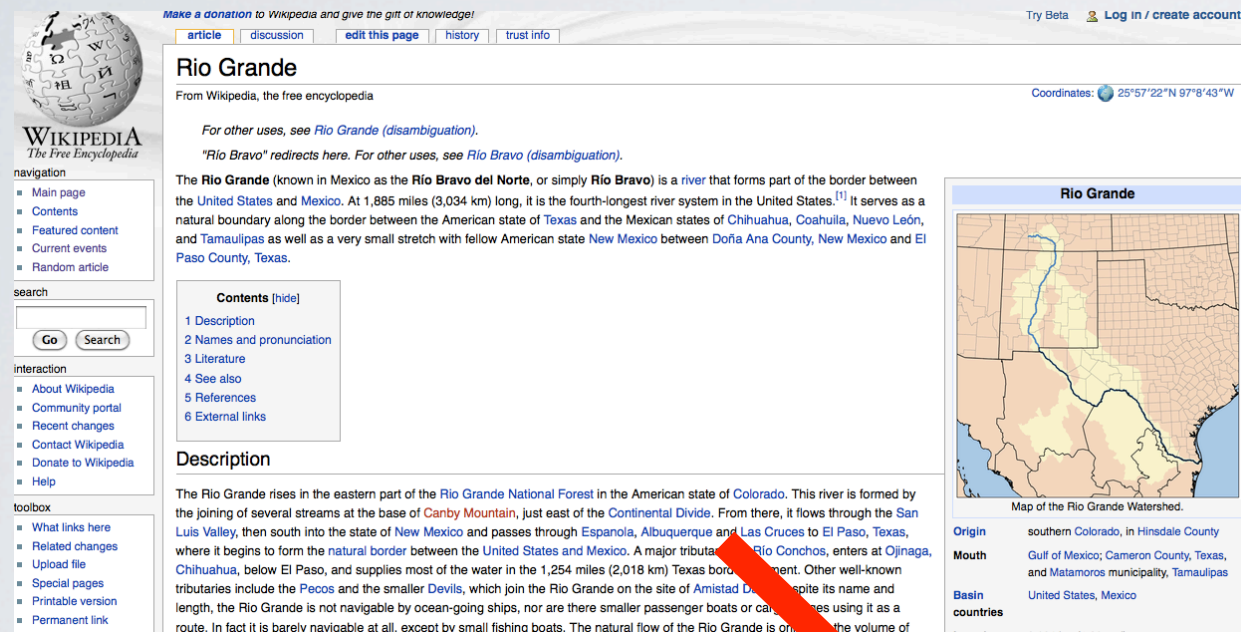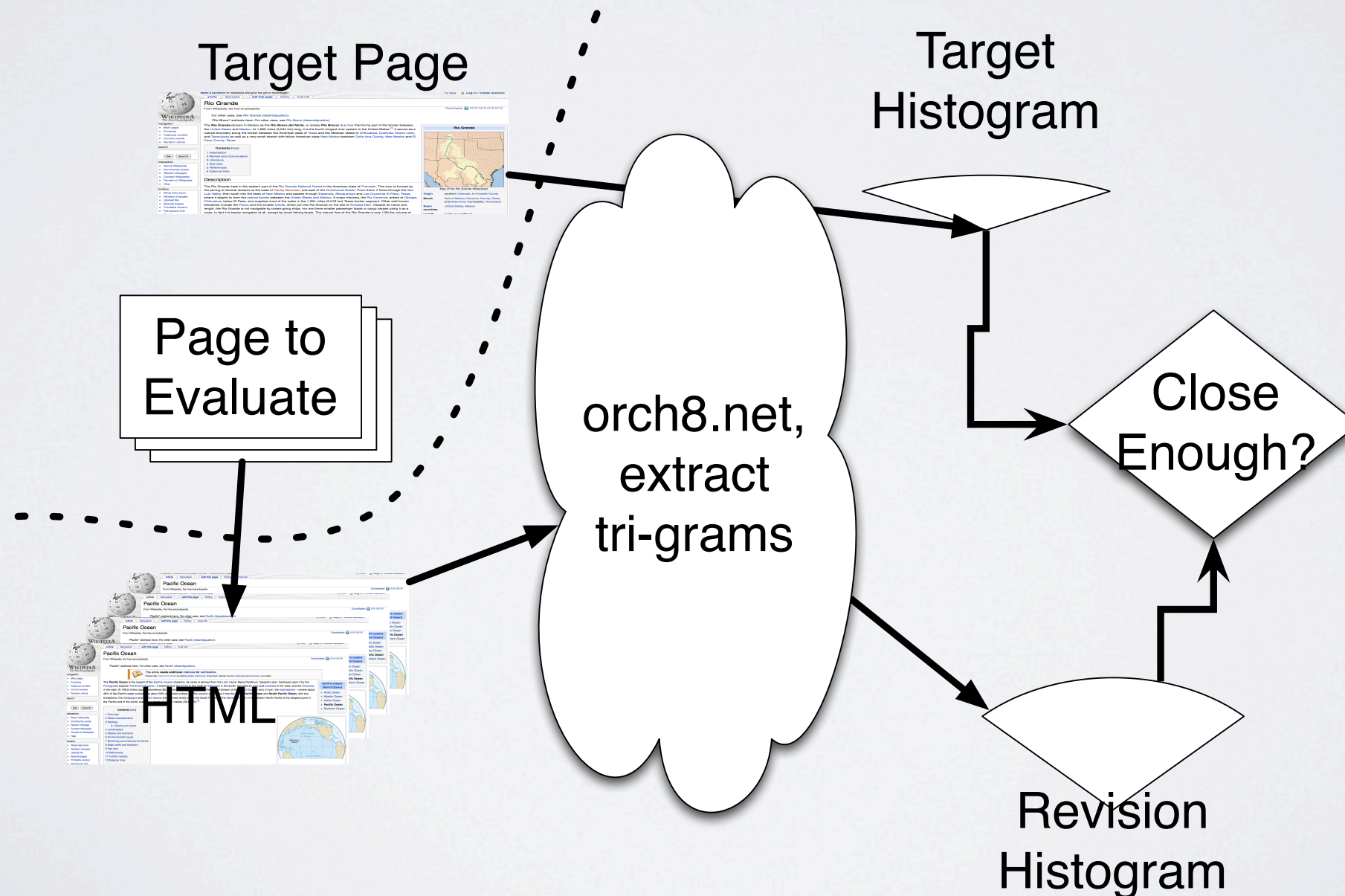## 2) From HTML to Semantic Digest



Rio Grande River
Texas State River
Big Bend Park

# OUR SOLUTION

1) Download and Render to HTML
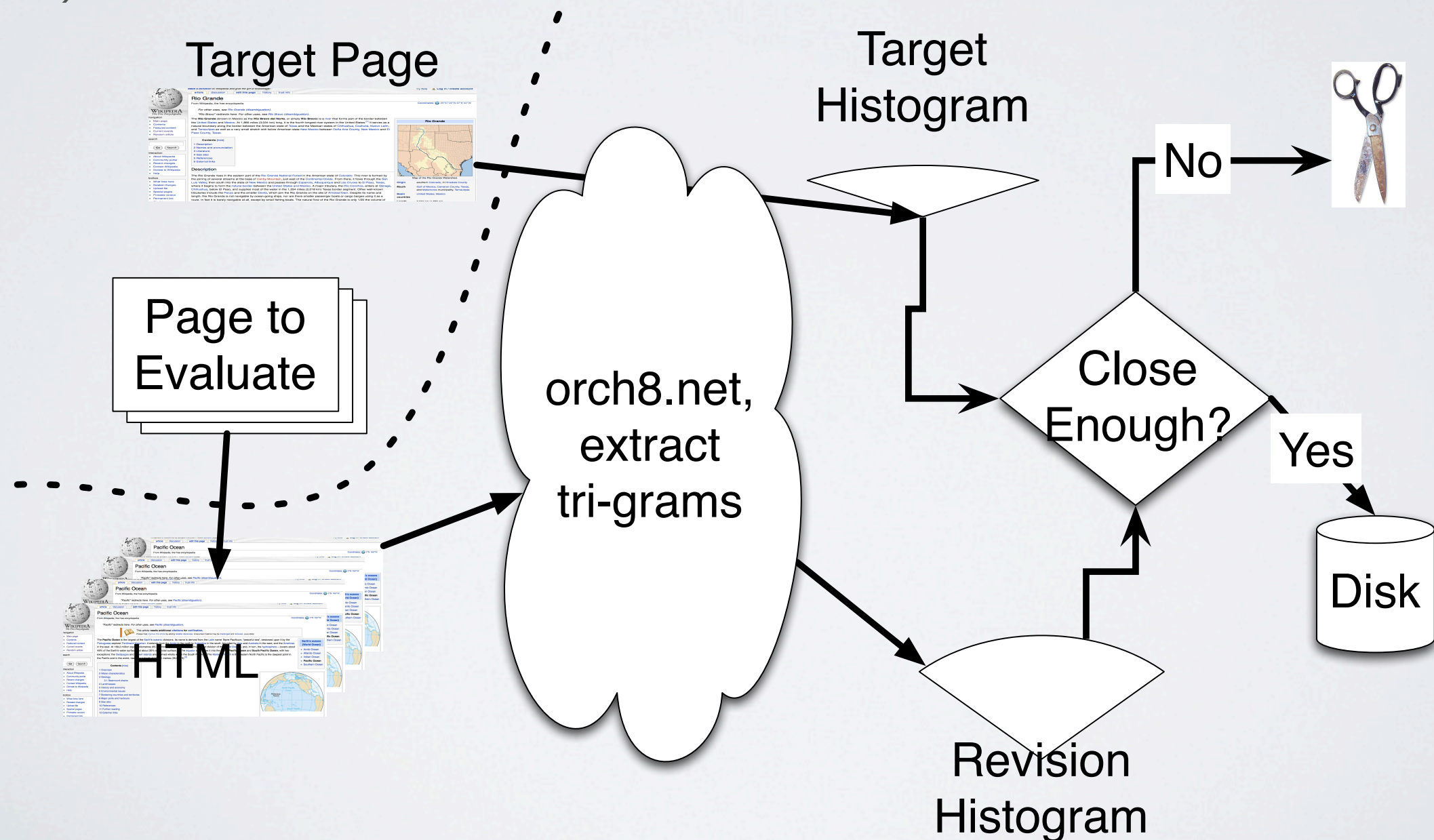2) From HTML to Semantic Digest
3) Euclidian Distance from the Target

**Target Page**

**Page to Evaluate**

**HTML**

**orch8.net, extract tri-grams**

**Target Histogram**

**Close Enough?**

**Revision Histogram**

# OUR SOLUTION

1) Download and Render to HTML
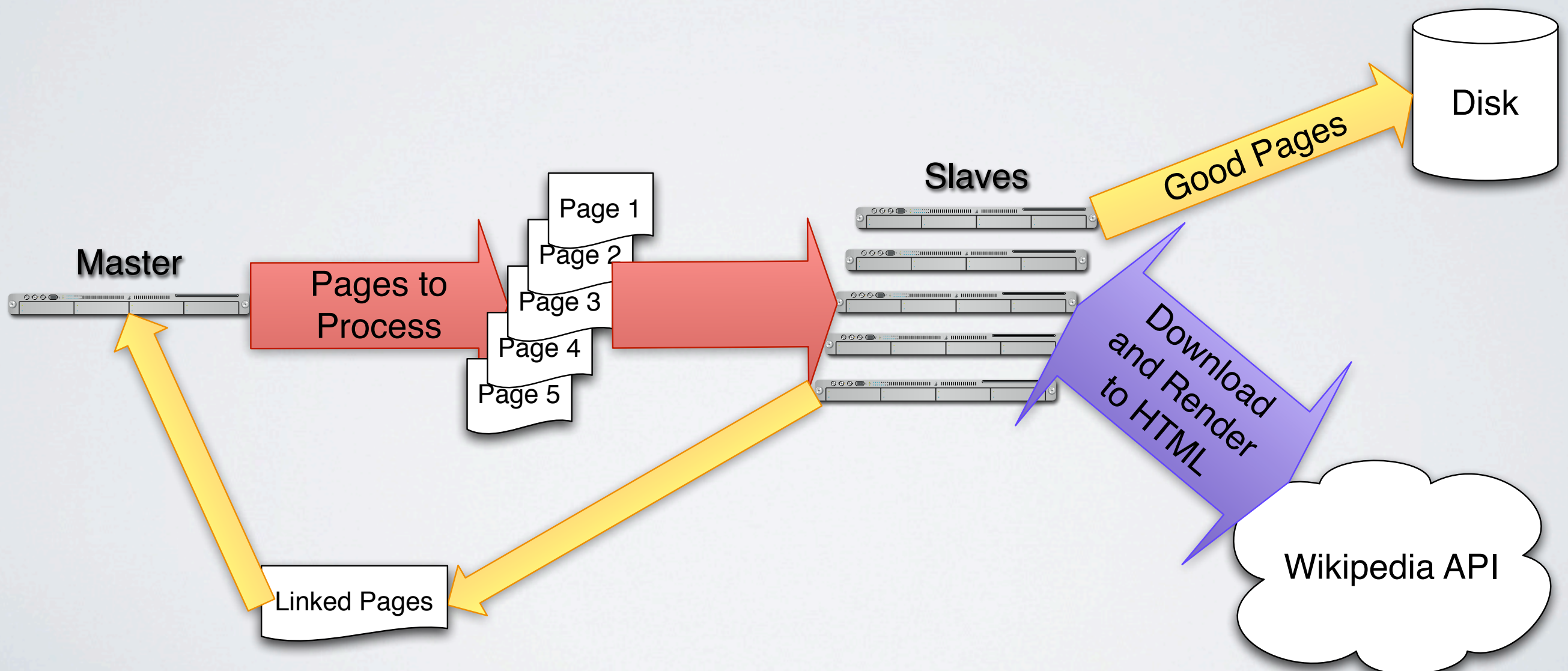2) From HTML to Semantic Digest
3) Euclidian Distance from the Target



432

# OUR SOLUTION

1) Download and Render to HTML
2) From HTML to Semantic Digest
3) Euclidian Distance from the Target
4) Extract Links and Recurse If Distance < Max

**Target Page**

**Page to Evaluate**

**HTML**

orch8.net, extract tri-grams

**Target Histogram**

**Revision Histogram**

No

Close Enough?

Yes

Disk

# DO THIS IN PARALLEL

OpenMPI allows us to run in a page-wise parallel fashion on a cluster.

# TESTING PLATFORM:
## CRAIGZCRUZER
## [WWW.SUPERTRIC.COM](WWW.SUPERTRIC.COM)
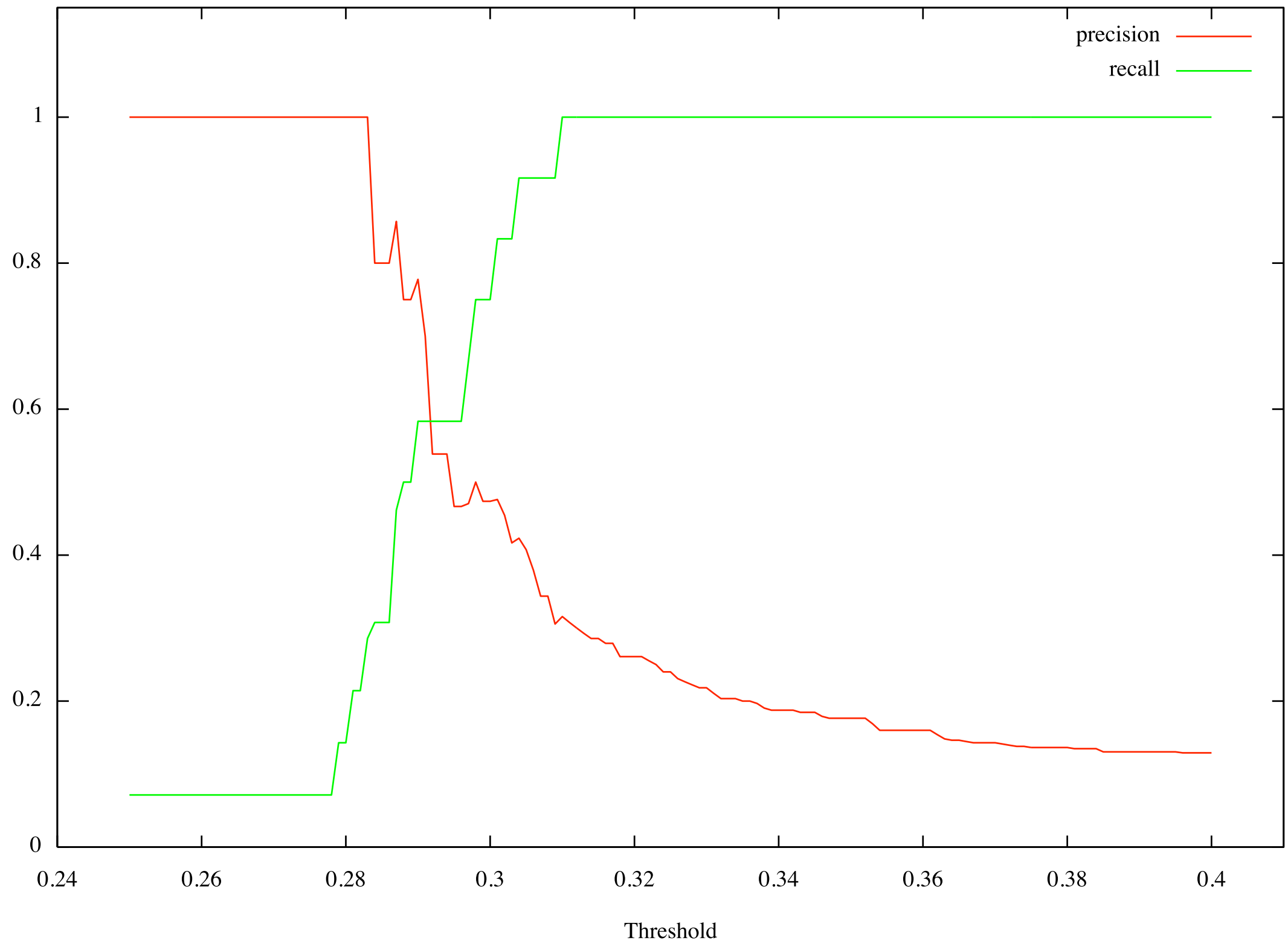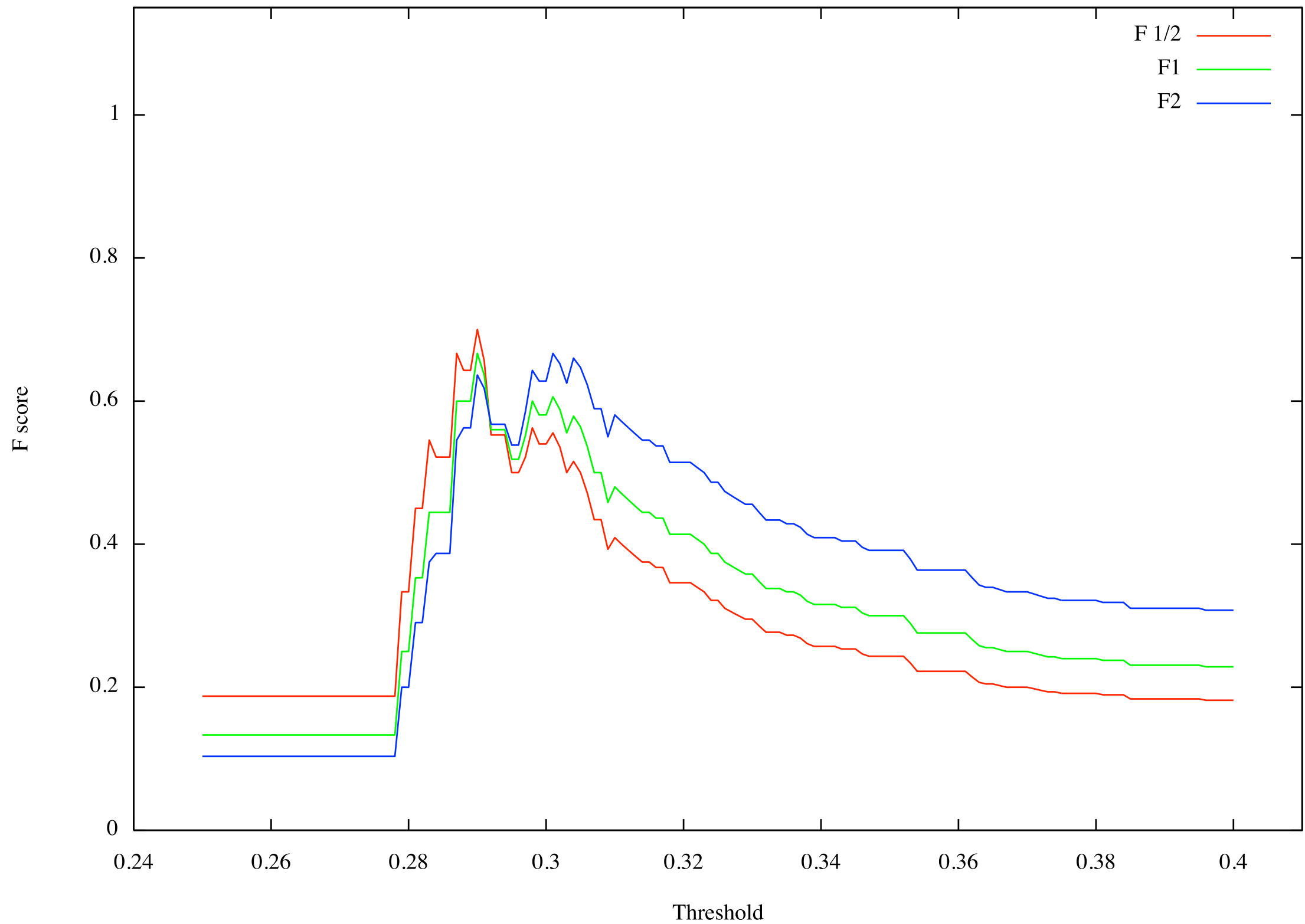
# HOW ACCURATE IS IT?
## METRICS (OR, A SHORT DIGRESSION INTO IR)

- Precision: the ratio of true positives to false positives

- Recall: the ratio of found positives to all positives

- $F_\beta$: How much to weight recall vs precision?

# RESULTS: PRECISION AND RECALL

# RESULTS: Fβ

# CONCLUSIONS

- SCrawl is a tool to extract what text MIGHT be interesting.

- Works though semantic meaning.

- Achieve high recall without sacrificing too much precision.

- Parallel processing for large datasets.